



DIGITAL PRESERVATION FOR INSTITUTIONAL REPOSITORIES

ELIZABETH LA BEAUD

DIGITAL LAB MANAGER, UNIVERSITY OF SOUTHERN MISSISSIPPI



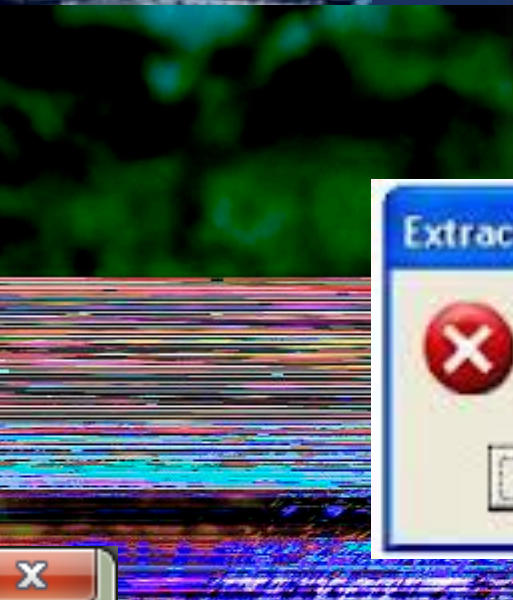
ACCESS VS. PRESERVATION

ACCESS

- relies on **cutting edge** technologies to provide best and fastest access at a point in time
- **selects** metadata needed to use and understand content
- access systems **deliver** objects with user-oriented services to make the objects
- purpose: **provide content to users**
- focus: **current users**

PRESERVATION

- relies upon **proven** technologies to preserve digital objects across generations of technology
- **accumulates** metadata over the life cycle to trace and preserve content
- preservation systems **create** new versions of digital objects for access to deliver as needs change over time
- purpose: **ensure long-term access**
- focus: **future** users



Microsoft Excel

The file is corrupt and cannot be opened.

OK

[Was this information helpful?](#)

Extraction Failed

File is corrupt

OK

RISKS TO DIGITAL CONTENT

- Obsolescence – hardware & software
- Loss of power
- Hardware failure
- Change and loss – accidental and intentional
- Bit rot
- User error
- Man-made and natural disasters
- Cyber attacks
- Format obsolescence
- Media failure
- Inappropriate access (i.e. confidential data)
- Non-compliance – standards and requirements
- Dust



PRESERVATION MAKES LONG-TERM ACCESS
POSSIBLE...



NDSA'S LEVELS OF DIGITAL PRESERVATION

Table 1: Version 1 of the Levels of Digital Preservation

	Level 1 (Protect your data)	Level 2 (Know your data)	Level 3 (Monitor your data)	Level 4 (Repair your data)
Storage and Geographic Location	<ul style="list-style-type: none"> - Two complete copies that are not collocated - For data on heterogeneous media (optical discs, hard drives, etc.) get the content off the medium and into your storage system 	<ul style="list-style-type: none"> - At least three complete copies - At least one copy in a different geographic location - Document your storage system(s) and storage media and what you need to use them 	<ul style="list-style-type: none"> - At least one copy in a geographic location with a different disaster threat - Obsolescence monitoring process for your storage system(s) and media 	<ul style="list-style-type: none"> - At least three copies in geographic locations with different disaster threats - Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems
File Fixity and Data Integrity	<ul style="list-style-type: none"> - Check file fixity on ingest if it has been provided with the content - Create fixity info if it wasn't provided with the content 	<ul style="list-style-type: none"> - Check fixity on all ingests - Use write-blockers when working with original media - Virus-check high risk content 	<ul style="list-style-type: none"> - Check fixity of content at fixed intervals - Maintain logs of fixity info; supply audit on demand - Ability to detect corrupt data - Virus-check all content 	<ul style="list-style-type: none"> - Check fixity of all content in response to specific events or activities - Ability to replace/repair corrupted data - Ensure no one person has write access to all copies
Information Security	<ul style="list-style-type: none"> - Identify who has read, write, move and delete authorization to individual files - Restrict who has those authorizations to individual files 	<ul style="list-style-type: none"> - Document access restrictions for content 	<ul style="list-style-type: none"> - Maintain logs of who performed what actions on files, including deletions and preservation actions 	<ul style="list-style-type: none"> - Perform audit of logs
Metadata	<ul style="list-style-type: none"> - Inventory of content and its storage location - Ensure backup and non-collocation of inventory 	<ul style="list-style-type: none"> - Store administrative metadata - Store transformative metadata and log events 	<ul style="list-style-type: none"> - Store standard technical and descriptive metadata 	<ul style="list-style-type: none"> - Store standard preservation metadata
File Formats	<ul style="list-style-type: none"> - When you can give input into the creation of digital files encourage use of a limited set of known open formats and codecs 	<ul style="list-style-type: none"> - Inventory of file formats in use 	<ul style="list-style-type: none"> - Monitor file format obsolescence issues 	<ul style="list-style-type: none"> - Perform format migrations, emulation and similar activities as needed

You have options:

- Buy, build, or join
- Local installation

- Open source
- Hosted solutions
- Proprietary



CONSIDER...

- Cost (available resources for preservation)
- Quantity (size and number of files)
- Expertise (skills required to manage)
- Partnerships (geographic distribution)
- Services (outsourcing)



WHERE IS YOUR DATA?





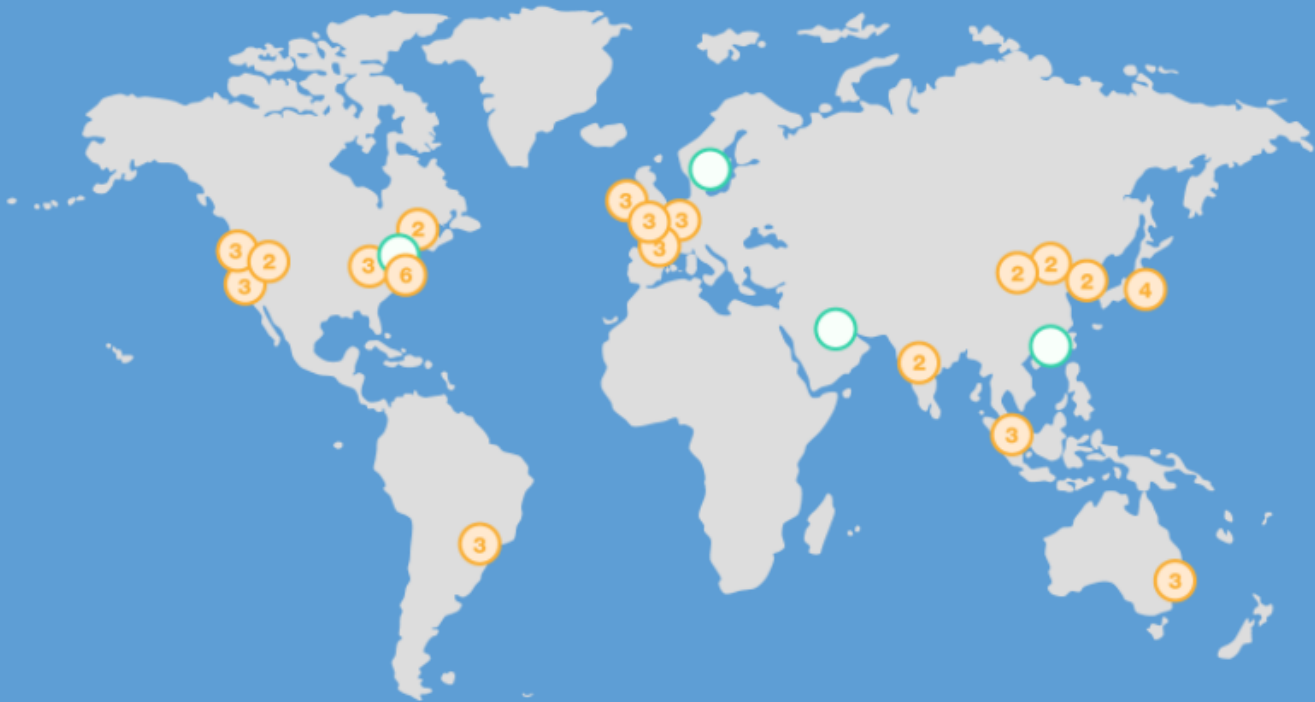
WHERE IS THE CLOUD?



Amazon Web Services

Global Network of Regions and Edge Locations

The AWS Cloud spans 52 Availability Zones within 18 geographic Regions around the world, with announced plans for 12 more Availability Zones and four more Regions in Bahrain, Hong Kong SAR, Sweden, and a second AWS GovCloud Region in the US.



Region & Number of Availability Zones

US East N. Virginia (6), Ohio (3)	Europe Frankfurt (3), Ireland (3), London (3), Paris (3)
US West N. California (3), Oregon (3)	South America São Paulo (3)
Asia Pacific Mumbai (2), Seoul (2), Singapore (3), Sydney (3), Tokyo (4)	AWS GovCloud (US-West) (2)
Canada Central (2)	
China Beijing (2), Ningxia (2)	



New Region (coming soon)

Bahrain
Hong Kong SAR, China
Sweden
AWS GovCloud (US-East)

AT LEAST TWO COPIES IN TWO LOCATIONS...

Storage and Geographic Location	<ul style="list-style-type: none">- Two complete copies that are not collocated- For data on heterogeneous media (optical discs, hard drives, etc.) get the content off the medium and into your storage system	<ul style="list-style-type: none">- At least three complete copies- At least one copy in a different geographic location- Document your storage system(s) and storage media and what you need to use them	<ul style="list-style-type: none">- At least one copy in a geographic location with a different disaster threat- Obsolescence monitoring process for your storage system(s) and media	<ul style="list-style-type: none">- At least three copies in geographic locations with different disaster threats- Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems
---------------------------------	--	---	--	---



PLAN FOR STORAGE SYSTEM REPLACEMENT





WHO HAS ACCESS TO YOUR DATA?





WHAT HAPPENS WHEN THEY LEAVE?



SECURITY

Information Security	<ul style="list-style-type: none">- Identify who has read, write, move and delete authorization to individual files- Restrict who has those authorizations to individual files	<ul style="list-style-type: none">- Document access restrictions for content	<ul style="list-style-type: none">- Maintain logs of who performed what actions on files, including deletions and preservation actions	<ul style="list-style-type: none">- Perform audit of logs
----------------------	---	--	--	---



Make sure

no one person

has write privileges
to all copies

FILE FIXITY & DATA INTEGRITY

File Fixity and Data Integrity	<ul style="list-style-type: none">- Check file fixity on ingest if it has been provided with the content- Create fixity info if it wasn't provided with the content	<ul style="list-style-type: none">- Check fixity on all ingests- Use write-blockers when working with original media- Virus-check high risk content	<ul style="list-style-type: none">- Check fixity of content at fixed intervals- Maintain logs of fixity info; supply audit on demand- Ability to detect corrupt data- Virus-check all content	<ul style="list-style-type: none">- Check fixity of all content in response to specific events or activities- Ability to replace/repair corrupted data- Ensure no one person has write access to all copies
--------------------------------	--	---	--	---

ALERTS YOU TO A PROBLEM

- Checksums
 - SHA256
 - MD5
- MD5 Summer
- AVPreserve's Fixity
- File Verifier ++
 - (these are windows examples)

```
fixity_2014-08-07-141939227000_cartoons - Notepad
File Edit Format View Help
Fixity report
Project name    cartoons
Algorithm used  md5
Date           2014-08-07
Total Files    7613
Confirmed Files 7613
Moved or Renamed Files 0
New Files      0
Changed Files  0
Removed Files  0
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0403.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0044.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1314.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1475.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1493.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1360.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0370.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1361.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0166.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0879.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0790.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1422.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0895.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0609.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1480.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0250.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1348.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0519.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0578.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0451.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1172.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0554.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0669.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1040.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0159.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0150.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1142.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0254.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0876.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0646.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0623.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0280.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0061.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0886.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0116.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC0995.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1046.tif
Confirmed File: T:\master tiffs\editorial cartoons\AAEC1312.tif
```




CAN YOUR VENDOR PROVE INTEGRITY ON DEMAND?

HOW OFTEN DO THEY CHECK FILE INTEGRITY?



FILE FORMAT SUSTAINABILITY

File Formats	<ul style="list-style-type: none">- When you can give input into the creation of digital files encourage use of a limited set of known open formats and codecs	<ul style="list-style-type: none">- Inventory of file formats in use	<ul style="list-style-type: none">- Monitor file format obsolescence issues	<ul style="list-style-type: none">- Perform format migrations, emulation and similar activities as needed
--------------	--	--	---	---



WHAT ARE YOU ACCEPTING?

HOW ARE YOU GOING TO MIGRATE IT OVER TIME?



TIERED OPTIONS

- **Tier 1**

- Open Sustainable File Formats (e.g., .txt)
- We will preserve your data to the best of our ability, and migrate the content over time to ensure accessibility. We'll keep a copy of your original file as well, just in case you want to revisit the original bitstream (1's and 0's).

- **Tier 2**

- Some well-known, but proprietary files (e.g., Microsoft Word)
- We will preserve your original file, but are limited in how we can maintain it. We'll always keep the 1's and 0's in the right order, but some formatting or other features may change.

- **Tier 3**

- Proprietary file formats with little to no open documentation (e.g., .spss)
- There are times when a specialized format is required for your data. In these cases, we will work with you. We will commit to preserving the original bitstream (1's and 0's) and the submitted metadata; however, while the files may or may not be renderable in the future, we can make no guarantees on the form or presentation it will take.

SUSTAINABLE FILE FORMAT CHARACTERISTICS

- Widely used by appropriate research community
- Non-proprietary
- Well documented, open standard
- Uncompressed (size permitting)
- Unencrypted
- Encoded in standard representation (e.g., ASCII, UTF-8)

METADATA

Metadata	<ul style="list-style-type: none">- Inventory of content and its storage location- Ensure backup and non-collocation of inventory	<ul style="list-style-type: none">- Store administrative metadata- Store transformative metadata and log events	<ul style="list-style-type: none">- Store standard technical and descriptive metadata	<ul style="list-style-type: none">- Store standard preservation metadata
----------	--	--	---	--



DESCRIBE IT SO YOU CAN FIND IT!



METADATA FIELDS

- Identifier
- Title
- Description
- Creator
- Department
- Legacy department
- Date
- coverage
- Subjects
- Geographic location
- Language
- Type
- File format
- Contributors
- Permissions
- Disciplines
- Related docs

RECORDS CHANGES OVER TIME

- Document metadata changes
- File format migrations
- IP problems & resolutions
- Deaccessions

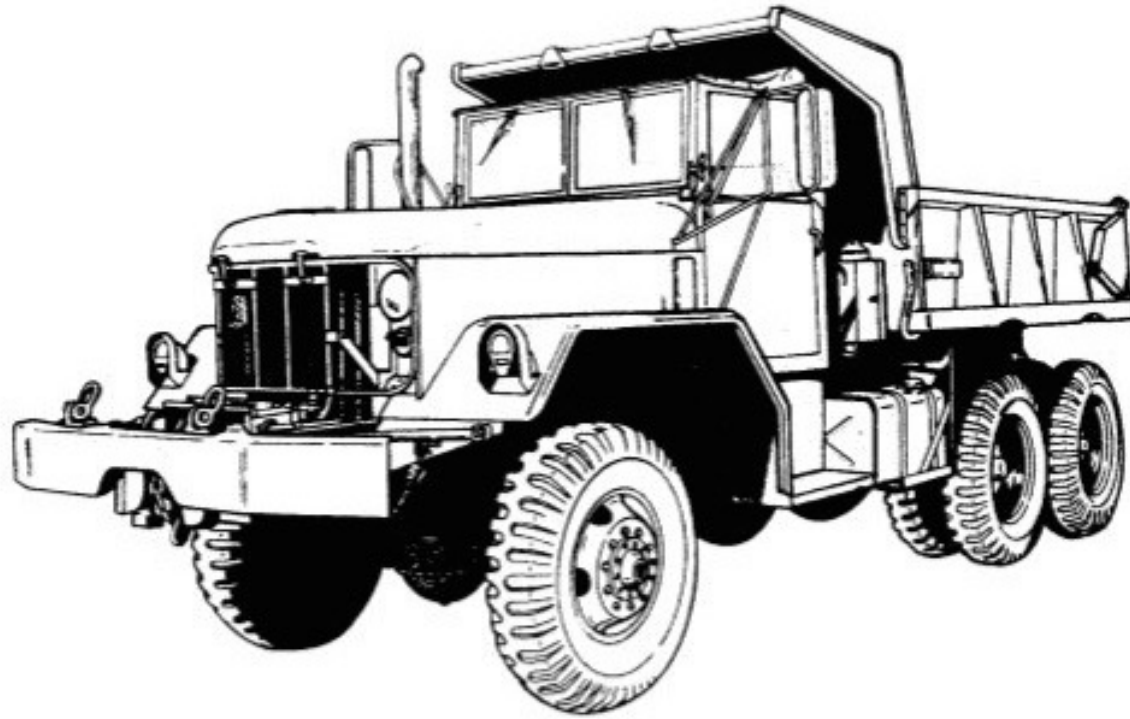


DOCUMENT, DOCUMENT, DOCUMENT...

POLICIES & PROCEDURES



TRUCK BOOK...



IT'S A COMMITMENT...



DPM Workshop



IT'S NOT A ONE AND DONE...





BUT...YOU DON'T HAVE TO DO IT ALONE...

RESOURCES





THANK YOU !!

ELIZABETH LA BEAUD

ELIZABETH.LABEAUD@USM.EDU